



**Prioritätsbescheinigung über die Einreichung
einer Patentanmeldung**

Aktenzeichen: 100 30 105.3
Anmeldetag: 19. Juni 2000
Anmelder/Inhaber: ROBERT BOSCH GMBH,
Stuttgart/DE
Bezeichnung: Spracherkennungseinrichtung
IPC: G 10 L 15/00

**Die angehefteten Stücke sind eine richtige und genaue Wiedergabe der ur-
sprünglichen Unterlagen dieser Patentanmeldung.**

München, den 15. März 2001
Deutsches Patent- und Markenamt
Der Präsident
Im Auftrag

Waasmaier

8L244504020

13. Juni 2000 - fle/ibr

Robert Bosch GmbH, 70442 Stuttgart

Spracherkennungseinrichtung

Die Erfindung bezieht sich auf eine Spracherkennungseinrichtung, bei der mindestens zwei Eingangssignale über jeweilige getrennte Kanäle parallel einer Erkennungseinrichtung mit einer Merkmalsextraktionseinrichtung zum Bilden von Merkmalvektoren, mit einer Transformationseinrichtung zum Bilden transformierter Merkmalvektoren und mit einem nachfolgenden Klassifikator zugeführt werden, der eine Klassifikation der zugeführten transformierten Merkmalvektoren trifft und den ermittelten Klassen entsprechende Ausgangssignale abgibt.

In modernen Systemen zur automatischen Spracherkennung wird oftmals durch eine lineare Transformation extrahierter Merkmale versucht, die Erkennungsleistung eines zugrundeliegenden Klassifikators zu verbessern. Die Transformation wird so gewählt, dass zum einen eine Dimensionsreduktion des Merkmalraumes stattfindet, aber zum anderen möglichst viel klassentrennende Information erhalten bleibt. Vielfach wird für diese Aufgabe die lineare Diskriminanzanalyse verwendet, wie z.B. in R. Haeb-Umbach, H. Ney,: Linear Discriminant Analysis for Improved Large Vocabulary Continuous Speech Recognition. In: Proceedings of the International Conference on Acoustics, Speech & Signal Processing

(ICASSP). I. 1992, S. 13-16; M. Finke, P. Geutner, H. Hild, T. Kemp, K. Ries, M. Westphal: The Karlsruhe-VerbMobil Speech Recognition Engine. In: Proceedings of the International Conference on Acoustics, Speech & Signal Processing (ICASSP). I. 1997, S. 83-86 sowie G. Ruske, R. Falthausen, T. Pfau: Extended Linear Discriminant Analysis (EL-DA) for Speech Recognition. In: Proceedings of the International Conference on Speech and Language Processing (ICSLP). 1998 näher ausgeführt ist.

Bekannt ist dabei eine Reduktion eines zusammengesetzten Merkmalvektors von meist 39 auf 32 Komponenten. Der ursprüngliche Merkmalvektor wird dabei aus der Kurzzeitleistung des Signals, 12 Mel-Frequenz-Cepstral-Koeffizienten (MFCC), wie in S. B. Davis, P. Mermelstein: Comparison of Parametric Representation for Monosyllabic Word Recognition in Continuously Spoken Sentences. IEEE Transactions on Acoustics, Speech, and Signal Processing ASSP-28 (1980), S. 357-366 angegeben, und deren erster und zweiter zeitlicher Ableitung gebildet. Die Merkmalsextraktion arbeitet in diesem Fall mit einem einzigen Eingangssignal. Typischerweise werden die Merkmale für Signalblöcke mit einer Länge von etwa 20 ms berechnet. Dies geschieht in einem reduzierten Zeittakt etwa alle 10 ms. Eine solche Verarbeitungskette ist in Fig. 3 dargestellt. Der Index k bezeichnet dabei einen hohen Zeittakt eines digitalisierten Sprachsignals, während der Index l einen reduzierten Zeittakt der Merkmalvektoren darstellt. Die nachfolgende Klassifikation verwendet zur Unterscheidung einzelner Klassen sogenannte verborgene Markov-Modelle oder einen Mustervergleich mit dynamischer Zeitanpassung. Auch werden künstliche neuronale Netze zur Klassifikation verwendet. Diese Klassifikatoren müssen anhand von Beispieldaten in einer Trainingsphase an die Klassifikationsaufgabe angepasst werden.

5 Stehen andererseits mehrere Eingangssignale zur Verfügung, werden diese meist
mittels eines Verfahrens zur mehrkanaligen Störgeräuschreduktion zu einem
einzigem störgeräuschreduzierten Signal zusammengeführt, so dass die Merkmal-
extraktionseinrichtung der Spracherkennungseinrichtung selbst wiederum ledig-
10 lich nur ein ihr zugeführtes Eingangssignal verarbeiten muss. Die verwendeten
Verfahren zur Störgeräuschreduktion nutzen dabei die Korrelation zwischen den
Signalen, wie in J. Allen, D. Berkley, J. Blauert: Multimicrophone signal-
processing technique to remove room reverberation from speech signals. Journal
of the Acoustical Society of America 62 (1977), Nr. 4, S. 912-915 und M. Dör-
15 becker, S. Ernst: Combination of Two-Channel Spectral Subtraction and Adap-
tive Wiener Post-Filtering for Noise Reduction and Dereverberation. In: Pro-
ceedings of EUSIPCO. 2. 1996, S. 995-998 angegeben, oder die Richtwirkung
sogenannter Mikrofon-Arrays aus, wie in M. Dörbecker: Small Microphone
Arrays with Optimized Directivity for Speech Enhancement. In: Proceedings of
20 the European Conference on Speech Communication and Technology (EURO-
SPEECH). 1. 1997, S. 327-330 und J. Bitzer, K. U. Simmer, K.-D. Kammeyer:
Multi-Microphone Noise Reduction Techniques for Hands-Free Speech Recog-
nition - A Comparative Study. In: Proceedings of the Workshop on Robust
Methods for Speech Recognition in Adverse Conditions. 1999, S. 171-174
25 angegeben. Diese Verfahren arbeiten entweder im Frequenzbereich mit etwa
128 bis 512 Frequenzbändern oder mit einer Filterung der Eingangssignale im
Zeitbereich. Diese Ansätze erfordern eine hohe Rechenleistung, insbesondere bei
Echtzeitimplementierung, da zur Verrechnung große Datenmengen anfallen. Die
Reduktion auf wenige Merkmale erfolgt erst nach der Kombination der Ein-
gangssignale.

5 Der Erfindung liegt die Aufgabe zugrunde, eine Spracherkennungseinrichtung bereitzustellen, die bei möglichst hoher Erkennungsrate einen möglichst geringen Aufwand hinsichtlich des Aufbaus und der Verarbeitungsleistung erfordert.

10 Diese Aufgabe wird mit den Merkmalen des Anspruchs 1 gelöst. Hiernach ist vorgesehen, dass die Merkmalsextraktionseinrichtung in den einzelnen Kanälen getrennt angeordnete Merkmalsextraktionsstufen aufweist, die mit ihren Ausgängen an die gemeinsame Transformationseinrichtung angeschlossen ist.

15 Mit diesem Aufbau der Spracherkennungseinrichtung und der damit verwirklichten Vorgehensweise werden die Eingangssignale in den einzelnen Kanälen direkt der Merkmalsextraktion unterzogen, wobei aus den Eingangssignalen möglichst viel Information für den Erkennungsprozess in den extrahierten Merkmalvektor einfließen soll. Die Kombination der Kanäle findet erst im Merkmalraum statt, wobei aus den Merkmalvektoren der einzelnen Kanäle ein einziger transformierter Merkmalvektor berechnet wird. Die Merkmalvektoren werden also unab-
20 hängig voneinander aus den Eingangssignalen berechnet und mittels einer Transformation zu einem gemeinsamen Merkmalvektor kombiniert.

25 Während des Betriebs der Spracherkennungseinrichtung erfolgt die Kombination der Merkmalvektoren durch eine einfache zeitinvariante Matrixoperation. Dies führt gegenüber dem bekannten adaptiven Verfahren der mehrkanaligen Störgeräuschreduktion zu einer deutlichen Reduzierung des Rechenaufwandes. Zum einen ist für das entwickelte Verfahren keine Adaption im laufenden Betrieb notwendig, zum anderen erfolgt die Reduktion auf wenige Merkmale und auf
30 einen reduzierten Zeittakt bereits vor der Kombination der Kanäle.

5 Es hat sich überraschend gezeigt, dass beim Trainieren der Spracherkennungseinrichtung unter den Bedingungen einer vorgesehenen Betriebssituation ohne Störgeräuschreduktion einerseits und dem Einsatz der Spracherkennungseinrichtung in der entsprechenden realen Situation ebenfalls ohne Störgeräuschreduktion andererseits sich sogar eine höhere Erkennungsrate ergibt als mit Störgeräuschreduktion beim Training und realen Einsatz. Soll aus irgendwelchen
10 Gründen doch eine Störgeräuschreduktion beim Training und beim realen Einsatz vorgenommen werden, so kann diese ohne großen Zusatzaufwand vor der Merkmalsextraktion in einzelnen Kanälen, also pro Kanal relativ einfach vorgenommen werden.

15 Eine vorteilhafte Ausgestaltung der Spracherkennungseinrichtung besteht darin, dass die Transformationseinrichtung eine lineare Transformationseinrichtung ist. Geeignete Maßnahmen bestehen dabei darin, dass die Transformationseinrichtung zum Durchführen einer linearen Diskriminanzanalyse (LDA) oder einer
20 Karhunen-Loève-Transformation ausgebildet ist.

Diese Wahl der Transformationseinrichtung bei der Entwicklung des Spracherkenners führt dazu, dass möglichst viel Information zur Unterscheidung der verschiedenen Klassen erhalten bleibt. Beim Einsatz der linearen Diskriminanzanalyse oder der Karhunen-Loève-Transformation sind Beispieldaten für den Aufbau der Transformationseinrichtung notwendig. Günstig ist die Verwendung der-
25 selben Daten, die zum Aufbau des Klassifikators verwendet sind.

30 Es gibt auch Erweiterungen der LDA, die hier eingesetzt werden können. Außerdem ist es denkbar, nichtlineare Transformationseinrichtungen zu wählen (z.B.

5 sogenannte "neuronale Netze"). Gemeinsam ist diesen Verfahren, dass sie Beispieldaten für den Aufbau benötigen.

Die Erkennungsrate wird weiterhin dadurch unterstützt, dass der Klassifikator unter Bedingungen trainiert ist, die einer vorgesehenen Einsatzsituation entsprechen.

Die Erfindung wird nachfolgend unter Bezugnahme auf die Zeichnungen näher erläutert. Es zeigen:

15 Fig. 1 ein Blockschaltbild für eine zweikanalige Spracherkennungseinrichtung,

Fig. 2 ein Blockschaltbild für eine mehrkanalige Spracherkennungseinrichtung und

20 Fig. 3 eine einkanalige Spracherkennungseinrichtung nach dem Stand der Technik.

25 Fig. 1 zeigt ein Blockschaltbild einer entwickelten Spracherkennungseinrichtung bzw. eines entsprechenden Verfahrens in zweikanaliger Ausführung, das heißt mit zwei Eingangssignalen y_1 und y_2 . Merkmalvektoren O_1 und O_2 werden aus den Eingangssignalen y_1 und y_2 mittels bekannter Verfahren der Merkmalsextraktion, z.B. MFCC, getrennt pro Kanal gewonnen. Aus der Folge dieser Merkmal-

vektoren wird durch eine vorzugsweise lineare Operation eine neue Folge transformierter Merkmalvektoren gebildet nach der Beziehung

$$O^t(l) = T \cdot \begin{bmatrix} O_1(l) \\ O_2(l) \end{bmatrix} \quad (1)$$

Die Matrixoperation wird für jeden Signalblock in einem reduzierten Zeittakt l durchgeführt. Die Dimension der Matrix T wird entsprechend gewählt, um eine Dimensionsreduktion zu erbringen. Besitzen die Merkmalvektoren U_1 und U_2 jeweils n_1 bzw. n_2 Komponenten und soll der transformierte Merkmalvektor lediglich n_t Koeffizienten umfassen, so muss die Matrix T die Dimension n_t mal $(n_1 + n_2)$ besitzen. Ein typisches Zahlenbeispiel ist $n_1 = 39$, $n_2 = 39$ und $n_t = 32$. Dann hat die Transformationsmatrix T die Dimension $32 \cdot 78$, und die Transformation erbringt eine Dimensionsreduktion von insgesamt 78 Komponenten in den Merkmalvektoren O_1 und O_2 auf 32 Komponenten im transformierten Merkmalvektor O^t .

Die Transformationsmatrix T wird anhand von Beispieldaten so eingestellt, dass der transformierte Merkmalvektor O^t die maximale Information zur Unterscheidung der einzelnen Klassen besitzt. Dazu können die bekannten Verfahren der linearen Diskriminanzanalyse oder der Karhunen-Loève-Transformation verwendet werden. Zum Training des Klassifikators KL werden die transformierten Merkmalvektoren $O^t(l)$ verwendet.

- 5 Als Erweiterung des Verfahrens können, wie in Fig. 2 dargestellt, auch mehr als zwei Kanäle miteinander kombiniert werden. Gleichung 1 geht dann über in

$$O^t(l) = T \cdot \begin{bmatrix} O_1(l) \\ \vdots \\ O_N(l) \end{bmatrix} \quad (2)$$

Die Dimension der Transformationsmatrix ist dann $n_t \times \left(\sum_{i=1}^N n_i \right)$, wobei n_i die Anzahl der Komponenten in dem Merkmalvektor O_i angibt.


- 15 Die in den Fig. 1 und 2 angegebenen Blöcke ME1, ME2, ME_k der den jeweiligen Kanälen zugeordneten Merkmalsextraktionsstufen, die zusammen die Merkmalsextraktionseinrichtung bilden, müssen nicht notwendigerweise für alle Eingangssignale y_1 , y_2 bzw. y_N gleich sein. Als Alternative kommen z.B. Merkmale in Betracht, die auf der sogenannten linearen Prädiktion beruhen, die in der Sprachkodierung ebenfalls Anwendung findet.

13. Juni 2000 - fle/ibr

Robert Bosch GmbH, 70442 Stuttgart



Ansprüche

- 5
1. Spracherkennungseinrichtung, bei der mindestens zwei Eingangssignale $[y_1(k) \dots y_N(k)]$ über jeweilige getrennte Kanäle parallel einer Erkennungseinrichtung mit einer Merkmalsextraktionseinrichtung zum Bilden von Merkmalvektoren, mit einer Transformationseinrichtung zum Bilden transformierter Merkmalvektoren $[O^t(l)]$ und mit einem nachfolgenden Klassifikator (KL) zugeführt werden, der eine Klassifikation der zugeführten transformierten Merkmalvektoren $[O^t(l)]$ trifft und den ermittelten Klassen entsprechende Ausgangssignale abgibt, dadurch gekennzeichnet,
- 10  dass die Merkmalsextraktionseinrichtung in den einzelnen Kanälen getrennt angeordnete Merkmalsextraktionsstufen (ME1 ... MEk) aufweist, die mit ihren Ausgängen an die gemeinsame Transformationseinrichtung (T) angeschlossen sind.
- 15
2. Spracherkennungseinrichtung nach Anspruch 1, dadurch gekennzeichnet, dass die Transformationseinrichtung (T) eine lineare Transformationseinrichtung ist.

5

3. Spracherkennungseinrichtung nach Anspruch 2, dadurch gekennzeichnet, dass die Transformationseinrichtung (T) zum Durchführen einer linearen Diskriminanzanalyse oder einer Karhunen-Loève-Transformation ausgebildet ist.

10



4. Spracherkennungseinrichtung nach einem der vorhergehenden Ansprüche, dadurch gekennzeichnet, dass die Transformationseinrichtung (T) anhand von Beispieldaten gewählt ist.

15

5. Spracherkennungseinrichtung nach einem der vorhergehenden Ansprüche, dadurch gekennzeichnet, dass der Klassifikator (KL) unter Bedingungen trainiert ist, die einer vorgesehenen Einsatzsituation entsprechen.

20



6. Spracherkennungseinrichtung nach einem der vorhergehenden Ansprüche, dadurch gekennzeichnet, dass den Merkmalsextraktionsstufen (ME1 ... ME_k) jeweils zugeordnete Störgeräuschreduktionsstufen vorgeschaltet sind.

13. Juni 2000 - fle/ibr

Robert Bosch GmbH, 70442 Stuttgart

Verfahren zur mehrkanaligen Merkmalsextraktion für die Spracherkennung

Zusammenfassung

Die Erfindung bezieht sich auf eine Spracherkennungseinrichtung, bei der mindestens zwei Eingangssignale $[y_1(k) \dots y_N(k)]$ über jeweilige getrennte Kanäle parallel einer Erkennungseinrichtung mit einer Merkmalsextraktions-
einrichtung zum Bilden von Merkmalvektoren, mit einer Transformations-
einrichtung zum Bilden transformierter Merkmalvektoren $[O^t(l)]$ und mit
5 einem nachfolgenden Klassifikator (KL) zugeführt werden, der eine Klassifi-
kation der zugeführten transformierten Merkmalvektoren $[O^t(l)]$ trifft und den
ermittelten Klassen entsprechende Ausgangssignale abgibt. Eine hohe
Erkennungsrate bei relativ geringem Aufwand des Aufbaus und der
10 Verarbeitung werden dadurch erreicht, dass die Merkmalsextraktions-
einrichtung in den einzelnen Kanälen getrennt angeordnete Merkmalsextrak-
tionsstufen (ME1 ... MEk) aufweist, die mit ihren Ausgängen an die gemein-
same Transformationseinrichtung (T) angeschlossen sind (Fig. 1).

1/1

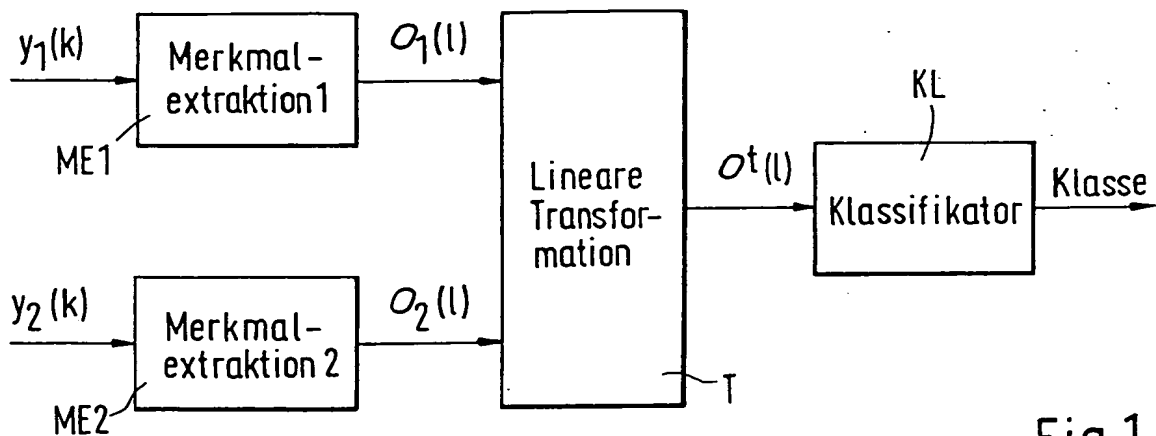


Fig.1

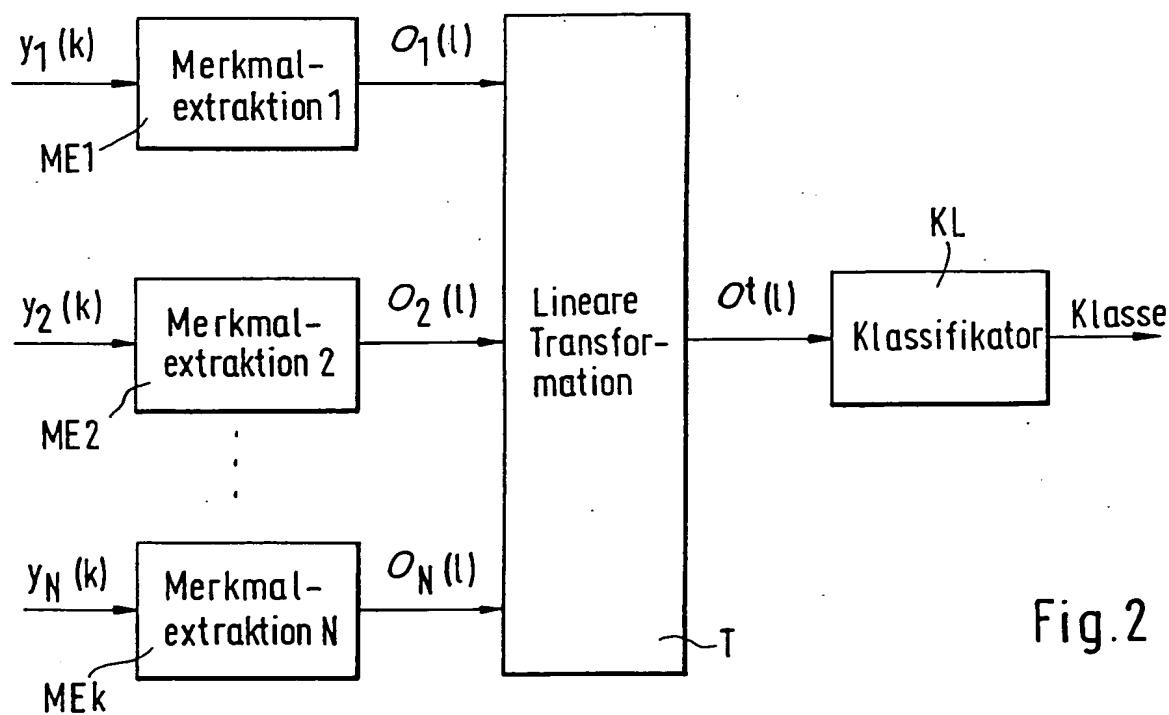


Fig.2

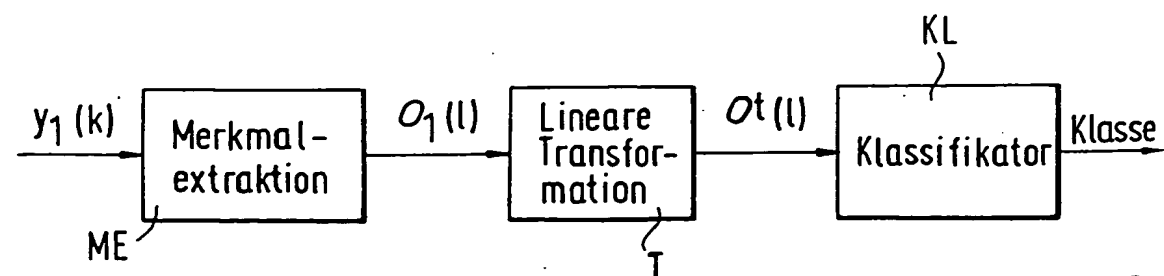


Fig.3